



大数据技术与应用专业规划教材

大数据技术与应用

◎ 娄岩 主编



清华大学出版社

大数据技术与应用专业规划教材

大数据技术与应用

娄 岩 主编

清华大学出版社
北 京

本书编委会

主 编：娄 岩

副 主 编：郑琳琳 徐东雨

编委成员(按姓氏笔画排列)：

丁 林 马 瑾 刘尚辉 张志常
李 静 庞东兴 曹 阳 霍 妍

前言	I
第 1 章 大数据概论	1
1.1 大数据技术简介	2
1.1.1 IT 产业的发展简史	2
1.1.2 大数据的主要来源	4
1.1.3 数据生成的 3 种主要方式	4
1.1.4 大数据的特点	5
1.1.5 大数据的处理流程	5
1.1.6 大数据的数据格式	6
1.1.7 大数据的基本特征	6
1.1.8 大数据的应用领域	7
1.2 大数据的技术架构	7
1.3 大数据的整体技术	8
1.4 大数据分析的 4 种典型工具简介	9
1.5 大数据未来发展趋势	10
1.5.1 数据资源化	10
1.5.2 数据科学和数据联盟的成立	10
1.5.3 大数据隐私和安全问题	11
1.5.4 开源软件成为推动大数据发展的动力	11
1.5.5 大数据在多方面改善人们的生活	12
本章小结	12
习题 1	12
第 2 章 大数据采集及预处理	14
2.1 数据采集简介	15
2.1.1 数据采集	15
2.1.2 数据采集的数据来源	15
2.1.3 数据采集的技术方法	17
2.2 大数据的预处理	18

2.3 大数据采集及预处理的主要工具·····	20
本章小结·····	29
习题2·····	29
第3章 大数据分析概论 ·····	31
3.1 大数据分析简介·····	32
3.1.1 大数据分析·····	32
3.1.2 大数据分析的基本方法·····	33
3.1.3 大数据处理流程·····	34
3.2 大数据分析的主要技术·····	36
3.2.1 深度学习·····	36
3.2.2 知识计算·····	37
3.3 大数据分析处理系统简介·····	39
3.3.1 批量数据及处理系统·····	39
3.3.2 流式数据及处理系统·····	40
3.3.3 交互式数据及处理系统·····	40
3.3.4 图数据及处理系统·····	40
3.4 大数据分析的应用·····	41
本章小结·····	43
习题3·····	43
第4章 大数据可视化 ·····	45
4.1 大数据可视化简介·····	45
4.2 大数据可视化工具 Tableau·····	50
本章小结·····	58
习题4·····	58
第5章 Hadoop 概论 ·····	59
5.1 Hadoop 简介·····	60
5.1.1 Hadoop 简史·····	60
5.1.2 Hadoop 应用和发展趋势·····	61
5.2 Hadoop 的架构与组成·····	62
5.2.1 Hadoop 架构介绍·····	63
5.2.2 Hadoop 组成模块·····	63
5.3 Hadoop 应用分析·····	65
本章小结·····	66
习题5·····	66

第 11 章 典型大数据解决方案	129
11.1 Intel 大数据	130
11.1.1 Intel 大数据解决方案	130
11.1.2 Intel 大数据相关案例	131
11.2 百度大数据	132
11.2.1 百度大数据引擎	132
11.2.2 百度大数据+平台	133
11.2.3 相关应用	133
11.2.4 百度预测的使用方法	135
11.3 腾讯大数据	137
11.3.1 腾讯大数据解决方案	137
11.3.2 相关实例	139
本章小结	140
习题 11	140
附录 A 习题答案	141
参考文献	151

第 1 章

大数据概论



导学

内容与要求

本章主要涉及大数据技术简介、大数据的技术架构、大数据的整体技术、大数据分析 4 种典型工具及大数据未来发展趋势,以便读者更好地了解什么是大数据技术。

“大数据技术简介”一节包含 IT 产业的发展简史、大数据的主要来源、数据生成的 3 种主要方式、大数据的特点、大数据的处理流程、大数据的数据格式、基本特征和应用领域。了解大数据的主要来源,掌握大数据的特点和大数据的处理流程。

“大数据的技术架构”一节介绍 4 层堆栈式技术架构,包括基础层、管理层、分析层和应用层。

“大数据的整体技术”一节介绍数据采集、数据存取、基础架构、数据处理、统计分析、数据挖掘、模型预测和结果呈现等大数据的整体技术。

“大数据分析的 4 种典型工具简介”一节介绍的工具包括 Hadoop、Spark、Storm 和 Apache Drill。

“大数据未来发展趋势”一节中简介数据资源化。随着大数据应用的发展,大数据资源成为重要的战略资源,数据成为新的战略制高点。

重点、难点

本章重点是了解大数据的特点、特征和大数据未来发展趋势,难点是了解大数据技术架构和整体技术。

7. Chukwa 是开源的_____,用于监控和分析大型分布式系统的数据。
8. Pig 的底层由一个编译器组成,它在运行的时候会产生一些_____程序序列。

二、简答题

1. 简述 Hadoop 第一代和第二代的区别。
2. 以表格形式阐述 HBase 与 Hive 的异同点。
3. 简述 Hadoop 在数据处理方面存在的问题。

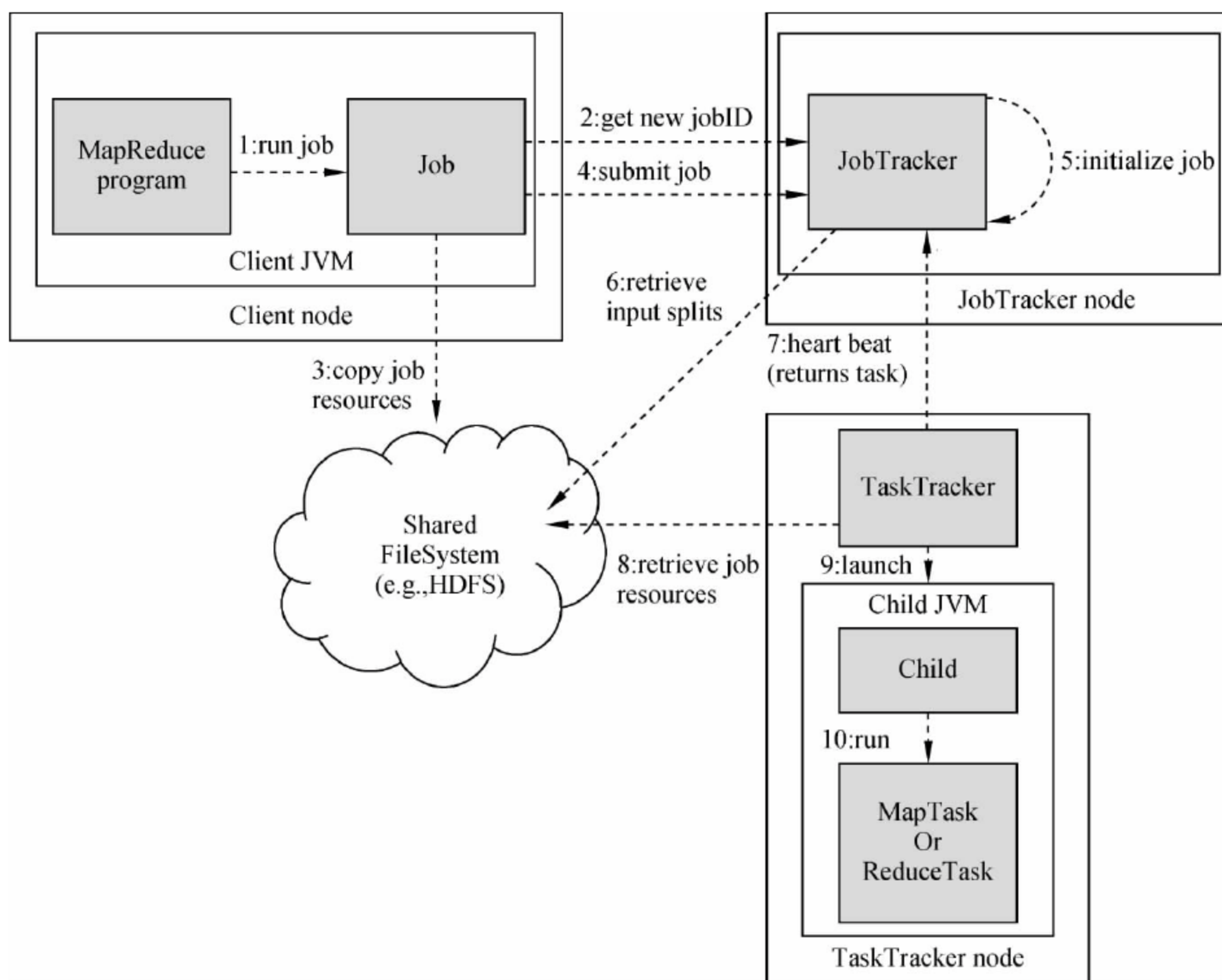


图 7-9 MapReduce 架构图

(3) TaskTracker: 将分配过来的数据片段执行 MapReduce 任务, 并保持与 JobTracker 通信。

(4) HDFS: 用来在其他节点间共享作业文件。

7.3.2 MapReduce 的工作流程

结合图 7-9, MapReduce 的工作流程可简单概括为以下 10 个工作步骤。

- (1) MapReduce 在客户端启动一个作业。
- (2) Client 向 JobTracker 请求一个 JobID。
- (3) Client 将需要执行的作业资源复制到 HDFS 上。
- (4) Client 将作业提交给 JobTracker。
- (5) JobTracker 在本地初始化作业。
- (6) JobTracker 从 HDFS 作业资源中获取作业输入的分割信息, 根据这些信息将作业分割成多个任务。
- (7) JobTracker 把多个任务分配给在与 JobTracker 心跳(即心跳信号)通信中请求任务的 TaskTracker。
- (8) TaskTracker 接收到新的任务之后会首先从 HDFS 上获取作业资源, 包括作业配置信息和本作业分片的输入。
- (9) TaskTracker 在本地登录子 JVM(Java Virtual Machine)。

(10) TaskTracker 启动一个 JVM 并执行任务,并将结果写回 HDFS。

本章小结

MapReduce 是 Hadoop 最重要的组成模块之一。MapReduce 由 Map 和 Reduce 两部分用户程序组成,利用框架在计算机集群上根据需求运行多个程序实例来处理各个子任务,然后再对结果进行归并输出。在实际的工作环境中,MapReduce 的分布式处理框架常用于分布式 Grep、分布式排序、Web 访问日志分析、反向索引构建、文档聚类、机器学习、数据分析、基于统计的机器翻译和生成整个搜索引擎的索引等大规模数据处理工作,并且已经在很多国内知名的互联网公司得到广泛地应用。

本章重点讲解了 MapReduce 的功能、技术特征、原理、架构和 workflows 等方面的知识。通过本章的学习,读者将会了解并掌握 MapReduce 的理论知识,为大数据方向的深入学习打下初步的基础。

习题 7

一、填空题

1. MapReduce 是_____。
2. MapReduce 适合处理各种类型的数据,包括_____数据。
3. MapReduce 采用了_____架构。
4. Map 功能是_____,然后返回一个基于这个处理的结果集。
5. Reduce 功能是_____,将多个 Map 的处理结果集进行分类和归纳。
6. _____意为分片,是 Map 任务最小的输入单位。
7. Shuffle 意为洗牌,一般包含本地化_____等操作。
8. 键值对是指 Key 和 Value 之间的_____关系,一个 Key 值对应一个 Value。
9. JobTracker 负责_____和 Tasks 的调度。
10. 用户提交的每一个 Job 会被划分成若干个_____。

二、简答题

1. 简述 MapReduce 的功能。
2. 简述 MapReduce 的技术特征。
3. 简述 MapReduce 的局限(至少列举出 3 点)。
4. 简述 Map 和 Reduce 的工作流程。
5. 简述 MapReduce 架构由哪些节点组成,各自的功能是什么。

第 8 章

NoSQL 概论



导学

内容与要求

本章主要介绍 NoSQL 的相关基础知识和 4 种类型数据管理方法(包括键值存储、列存储、面向文档存储和图形存储)的特点、数据管理的基本原理及典型工具。

在“NoSQL 简介”一节中,介绍 NoSQL 的含义、产生与特点。

在“NoSQL 技术基础”一节中,介绍一些与 NoSQL 相关的基本知识,包括一致性策略、分区与放置策略、复制与容错技术和缓存技术等。

在“NoSQL 的种类”一节中,介绍 NoSQL 的 4 种主要分类。

在“典型的 NoSQL 工具”一节中,针对 4 种类型的数据存储方式,分别介绍其典型工具。

重点、难点

本章重点是掌握 NoSQL 的基本知识以及分类,难点是 4 种不同类型的数据管理方法的工作原理及典型工具。

NoSQL 越来越多地被认为是关系型数据库的可行替代品,特别适用于大数据的存储。传统的关系型数据库因其对数据模式的约束程度高和对分布式存储的支持度差等因素,已经无法满足复杂、海量的数据存储。针对目前数据表现出的数量大、结构复杂、格式多样、存储要求不一致等特点,许多新兴的打破关系模型的数据存储方案应运而生,人们将其称为 NoSQL。通常情况下,人们把 NoSQL 解释为非结构化或是非关系型数据管理方法,其实更加准确的解释应该是 NoSQL--Not Only SQL,即不仅仅是关系型数据。

8.1 NoSQL 简介

8.1.1 NoSQL 的含义

NoSQL 泛指非关系型的数据管理技术。如果说 Hadoop 是一个产品,那么 NoSQL 就是一项技术。实际上,和处理常规数据一样,任何为处理大数据而服务的产品也都要选择符合实际情况的数据管理方式。由于网络上数据量激增,传统关系型数据库不能满足生活、生产需要,越来越多的人开始放弃严整、规矩的关系模型,另辟蹊径地去拓展研发新型的数据存储方式,如键值存储、列存储、面向文档存储和图形存储等,这些都属于 NoSQL 的范畴。

HDFS 在 Hadoop 中扮演数据存储的角色,可以将任何类型的文件按照分布式的方法进行存储。而 NoSQL 更侧重于数据管理层,可以应用于结构化、半结构化和非结构化数据存储。举一个例子,Hadoop 中的 HBase 正是采用 NoSQL 中的列存储方式对数据进行管理的。在 Hadoop 的架构中,Hbase 利用 HDFS 文件系统中存放的数据来解决特定的数据处理问题。这期间,HDFS 为 HBase 提供了高可靠性的底层存储支持,MapReduce 为 HBase 提供了高性能的计算能力。

8.1.2 NoSQL 的产生

随着大数据时代的到来及互联网 Web 2.0 网站的兴起,传统的关系型数据库在应付海量数据存储和读取,以及超大规模、高并发的 Web 2.0 纯动态网站的数据处理方面已经显得力不从心,同时也暴露出很多难以克服的问题。而非关系型的数据管理方法则由于其本身的特点得到了非常迅速的发展。NoSQL 技术的产生就是为了应对这一挑战。NoSQL 的概念最初在 2009 年被提出,对传统的数据管理方式是一次颠覆性的改变。

NoSQL 有很多种存储方式,拥有很多家族成员,NoSQL 的中文网站如图 8-1 所示,其中包括键值存储、面向文档存储、列存储、图形存储和 xml 数据存储等。其实在 NoSQL 的概念被提出之前,这些数据存储方式就已经被用于各种系统当中,只是很少被用于 Web 互联网应用中。

NoSQL 兴起的主要原因主要是传统的关系型数据库在网络数据存取上遇到了瓶颈。不得不说,传统的关系型数据库具有卓越的性能,高稳定性,且使用简单,功能强大,这使得传统的关系型数据库在 20 世纪 90 年代,网站访问数据量不是很大的情况下,发挥了令人瞩目的作用。

面临这些大数据管理的困扰,非关系型数据管理方式越来越被人们重视,并迅速发展。人们把这些有别于传统关系型数据库的数据管理技术统称为 NoSQL 技术。

在这里可以看到 NoSQL 的多个种类及各自的典型产品。

8.1.3 NoSQL 的特点

NoSQL 技术之所以能够在大数据冲击互联网的情况下脱颖而出,主要是因为其具有以下特点。



图 8-1 NoSQL 中文网站

(1) 易扩展性。尽管 NoSQL 数据库种类繁多,但是它们都有一个共同的特点,就是没有了关系型数据库中的数据与数据之间的关系。很显然,当数据之间不存在关系时,数据的可扩展性就变得可行。

(2) 数据量大,性能高。NoSQL 数据库都具有非常高的读写性能,尤其在大数据量下,同样表现优秀。这得益于它的无关系性,数据之间的结构简单。一般情况下,关系型数据库使用的是 Cache 在“表”这一层面的更新,是一种大粒度的 Cache 更新,当网络上的数据发生频繁交互时,就表现出了明显劣势。而 NoSQL 使用的是 Cache 在“记录”层面的更新,是一种细粒度的 Cache 更新,所以 NoSQL 在这个方面上也显示了较高的性能特点。

(3) 灵活的数据模型。由于 NoSQL 无须事先为要存储的数据建立字段,所以在应用中随时可以存储自定义的数据格式。而在关系数据库里,增删字段是一件非常麻烦的事情,尤其对数据量非常大的表而言,随时更改表结构几乎是无法实现的。而这一点在大数据量的 Web 2.0 时代尤为重要。

(4) 高可用性。NoSQL 在不太影响性能的情况,就可以方便地实现高可用的架构,比如 Cassandra、HBase 模型等。

8.2 NoSQL 技术基础

NoSQL 技术对大数据的管理是怎么实现的呢? 其中又要遵循哪些基本原则呢? 本节为读者在大数据的一致性策略、大数据的分区与放置策略、大数据的复制与容错技术及大数

据的缓存技术等方面进行介绍。

8.2.1 大数据的一致性策略

在大数据管理的众多方面,数据的一致性理论是实现海量数据进行管理的最基本的理论。学习这部分内容有利于读者对本章内容的阅读和深化理解。

分布式系统的 CAP 理论是构建 NoSQL 数据管理的基石。CAP,即一致性(Consistency)、可用性(Availability)和分区容错性(Partition Tolerance),如图 8-2 所示。

1. 一致性

一致性是指在分布式系统中的所有数据备份,在同一时刻均为同样的值。也就是当数据执行更新操作时,要保证系统内的所有用户读取到的数据是相同的。

2. 可用性

可用性是指在系统中任何用户的每一个操作均能在一定的时间内返回结果,即便当集群中的部分节点发生故障时,集群整体仍能响应客户端的读写请求。这里要强调“在一定时间内”,而不是让用户遥遥无期地等待。

3. 分区容错性

以实际效果而言,分区相当于对通信的时限要求。系统如果不能在时限内达成数据一致性,就意味着发生了分区的情况,必须就当前操作在一致性和可用性之间做出选择。

从上面的解释不难看出,系统不能同时满足一致性、可用性和分区容错性这 3 个特性,在同一时间只能满足其中的两个,如图 8-3 所示。因此系统设计者必须在这 3 个特性中做出抉择。



图 8-2 CAP 理论 3 个特性

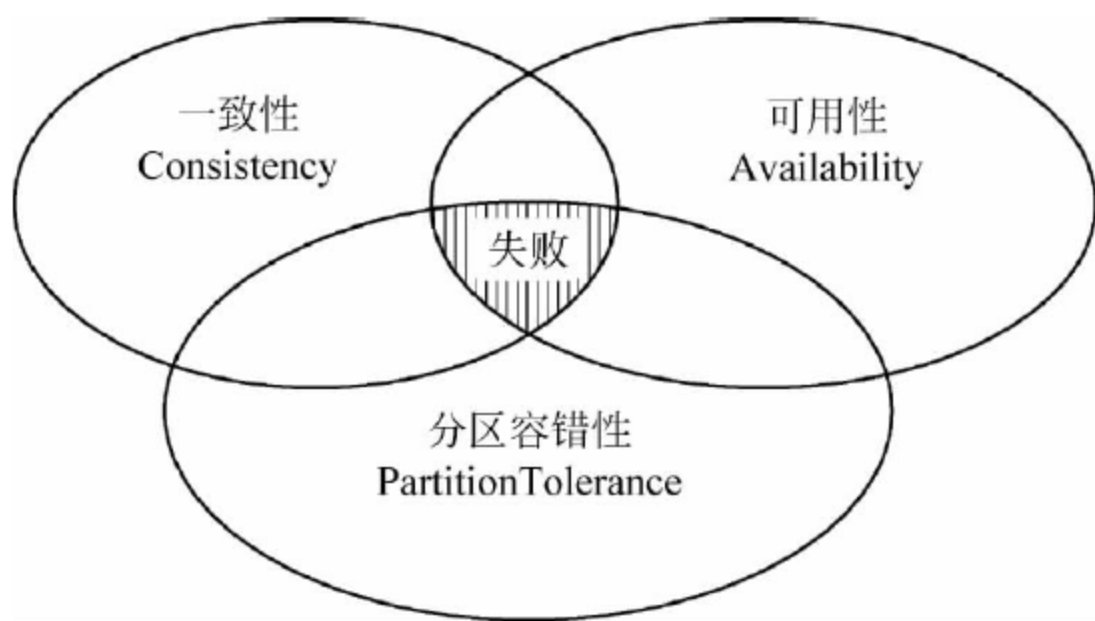


图 8-3 CAP 理论 3 个特性之间的关系

8.2.2 大数据的分区与放置策略

在大数据时代,如何有效地存储和处理海量的数据显得尤为重要。如果使用传统方法

处理这些数据,所消耗的时间代价将十分巨大,这是人们无法接受的,所以必须打破传统的将所有数据都存放在一处,每次查找、修改数据都必须遍历整个数据集的方法。数据分区技术与放置策略的出现正是为了解决数据存储空间不足及如何提高数据库性能等方面问题的。

1. 大数据分区技术

通俗地讲,数据分区其实就是“化整为零”,通过一定的规则将超大型的数据表分割成若干小块来分别处理。表进行分区时需要使用分区键来标志每一行属于哪一个分区,分区键以列的形式保存在表中。

数据分区可以提高数据的可管理性,改善数据库性能和数据可用性,缩小了每次数据查询的范围,并且在对数据进行维护时,可以只针对某一特定分区,大幅地提高数据维护的效率。

下面介绍几种常见的数据分区算法。

1) 范围分区

范围分区是最早出现的数据分区算法,也是最为经典的一个。所谓范围分区,就是将数据表内的记录按照某个属性的取值范围进行分区。

2) 列表分区

列表分区主要应用于各记录的某一属性上的取值为一组离散数值的情况,且数据集合中该属性在这些离散数值上的取值重复率很高。采用列表分区时,可以通过所要操作的数据直接查找到其所在分区。

3) 哈希分区

哈希分区需要借助哈希函数,首先把分区进行编号,然后通过哈希函数来计算确定分区内存储的数据。这种方法要求数据在分区上的分布是均匀的。

以上3种分区算法的特点和适用范围各异,在选择使用时应充分地考虑实际需求和数据表的特点,这样才能真正发挥数据分区在提高系统性能上的作用。

2. 大数据放置策略

为解决海量数据的放置问题,涌现了很多数据放置的算法,大体上可以分为两大类:顺序放置策略和随机放置策略。采用顺序放置策略是将各个存储节点看成是逻辑有序的,在对数据副本进行分配时先将同一数据的所有副本编号,然后采用一定的映射方式将各个副本放置到对应序号的节点上;随机放置策略通常是基于某一哈希函数来实现对数据的放置的,所以这里所谓的随机其实也是有规律的,很多时候称其为伪随机放置策略。

8.2.3 大数据的复制与容错技术

在大数据时代,每天都产生需要处理的大量的数据,在处理数据的过程中,难免会有差错,这可能会导致数据的改变和丢失。为了避免这些数据错误的出现,必须对数据进行及时的备份,这就是数据复制的重要性。同时,一旦出现数据错误,系统还要具备故障发现及处理故障的能力。

数据复制技术在处理海量数据过程中虽然是必不可少的,但是,对数据进行备份也要付出相应的代价。首先,数据的备份带来了大量的时间代价和空间代价;其次,为了减少时间

和空间上的代价,研究人员投入大量的时间、人力和物力来研发提升新的数据复制策略;另外,在数据备份的过程中往往会出现意想不到的差错,此时就需要数据容错技术和相应的故障处理方案进行辅助。

构成分布式系统的计算机五花八门,每台计算机又是由各式各样的软硬件组成的,所以在整个系统中可能随时会出现故障或错误。这些故障和错误往往是随机产生的,用户无法做到提前预知,甚至是当问题发生时都无法及时察觉。如果一个系统能够对无法预期的软硬件故障做出适当的对策和应变措施,那么就可以说这个系统具备一定的容错能力。

系统故障主要可以分为以下几类,如表 8-1 所示。

表 8-1 分布式环境下的系统故障类型

故障类型	故障子类	故障语义
崩溃故障	失忆型崩溃	服务器崩溃(停机),但停机前工作正常
		服务器只能从初始状态,遗忘了崩溃前的状态
	中顿型崩溃	服务器可以从崩溃前的状态启动
	停机型崩溃	服务器完全停机
失职故障	接收型失职	服务器对输入的请求没有响应
		服务器无法接收信件
	发送型失职	服务器无法发送信件
应答故障	返回值故障	服务器对服务请求做出错误反应
		返回值出现错误
	状态变迁故障	服务器偏离正确的运行轨迹
时序故障		服务器反应迟缓,超出规定的时间间隔
随意故障		服务器在任意时间产生的随意错误

处理故障的基本方法有主动复制、被动复制和半主动复制。所谓主动复制指的是所有的复制模块协同进行,并且状态紧密同步。被动复制是指只有一个模块为动态模块,其他模块的交互状态由这一模块的检查单定期更新。半主动复制是前两种的混合方法,所需的恢复开销相对较低。

8.2.4 大数据的缓存技术

单机的数据库系统引入缓存技术是为了在用户和数据库之间建立一层缓存机制,把经常访问的数据常驻于内存缓冲区,利用内存高速读取的特点来提高用户对数据查询的效率。在分布式环境下,由于组成系统的各个节点配置和使用的数据库系统及文件系统不尽相同,要想在这样复杂的环境下提高对海量数据的查询效率,仅仅依靠单机的缓存技术就行不通了。

与单机的缓存技术目的相同,分布式缓存技术的出现也是为了提高系统的数据查询性能。另外,为整个系统建立一层缓冲,也便于在不同节点之间进行数据交换。分布式缓存可以横跨多个服务器,所以可以灵活地进行扩展。

从图 8-4 中不难看出,如果各种.NET 应用、Web 服务和网格计算等应用程序在短时间内集中频繁的访问数据库服务器,很有可能会导致其瘫痪而无法工作。如果在应用程序和数据库之间加上一道缓冲屏障则可以解决这一问题。

时数据库中的数量也会急剧上升。

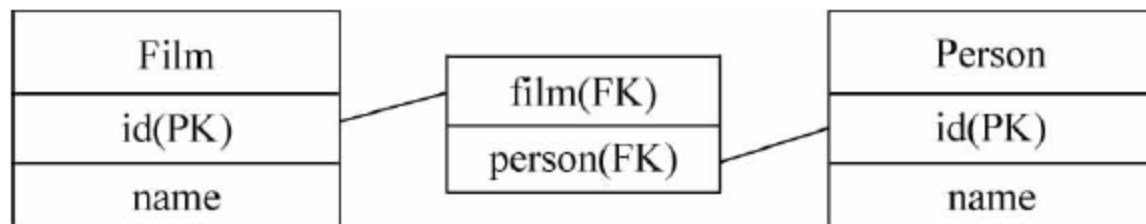


图 8-6 关系模型中的表及表间联系

除了性能之外,表的数量也是一个非常让人头疼的问题。刚刚仅仅是举了一个具有 4 个实体的例子——人、电影、电视剧和影视公司,现实生活中的例子可不是这么简单。不难看出,当需要描述大量关系时,传统的关系型数据库显得不堪重负,它更擅长的是实体较多但关系简单的情况。而对于一些实体间关系比较复杂的情况,高度支持关系的图形存储才是正确的选择。它不仅仅可以为人们带来运行性能的提升,更可以大大地提高系统开发效率,减少维护成本。

在需要表示多对多关系时,常常需要创建一个关联表来记录不同实体的多对多关系,而且这些关联表常常不用来记录信息。如果两个实体之间拥有多种关系,那么就需要在它们之间创建多个关联表。而在一个图形数据库中,只需要标明两者之间存在着不同的关系,例如,用 DirectedBy 关系指向电影的导演,或用 ActBy 关系来指定参与电影拍摄的各个演员,同时在 ActBy 关系中,更可以通过关系中的属性来表示其是否是该电影的主演。而且从上面所展示的关系的名称上可以看出,关系是有向的。如果希望在两个节点集间建立双向关系,就需要为每个方向定义一个关系。这两者的比较如图 8-7 所示。

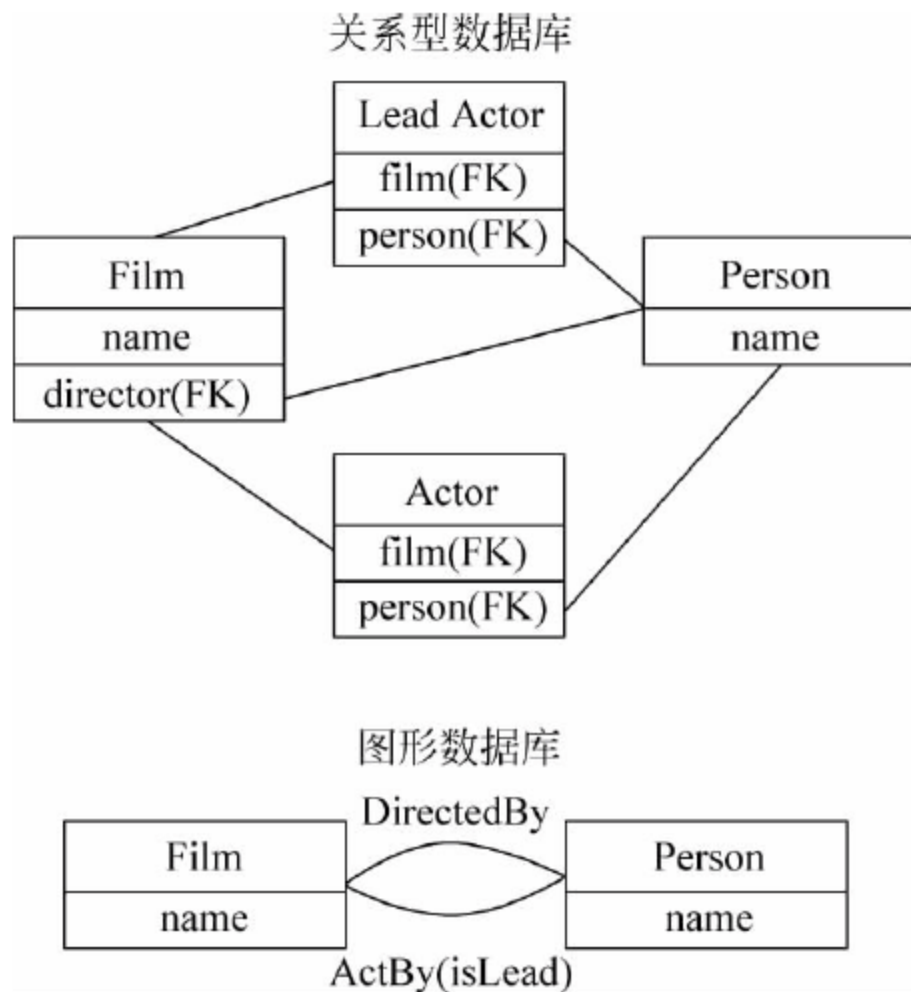


图 8-7 关系模型与图形存储的比较

8.4 典型的 NoSQL 工具

由于大数据时代刚刚到来,基于各类数据模型开发的数据库系统层出不穷,各个公司机构之间的竞争十分激烈。这一节将介绍目前实际应用中比较典型的 3 个 NoSQL 工具,以

此来代表 4 种不同的 NoSQL 数据管理类型。

8.4.1 Redis

Redis 是一个典型的开源 Key-Value 数据库。目前 Redis 的最新版本为 3.2.0,如图 8-8 所示。用户可以在 Redis 官网“<http://redis.io/download>”上获取最新的版本代码。



图 8-8 Redis 使用界面

1. Redis 的运行平台

Redis 可以在 Linux 和 Mac OS X 等操作系统下运行使用,其中 Linux 为主要推荐的操作系统。虽然官方没有提供支持 Windows 的版本,但是微软开发并维护一个 Win-64 的 Redis 端口。

2. Redis 的特点

(1) 支持存储的类型多样。与传统的关系型数据库或是其他非关系型数据库相比,Redis 支持存储的 Value 类型是非常多样的,不限于字符串,还包括 String(字符串)、Hash(哈希)、List(链表)、Set(集合)和 Zset(有序集合)等。

(2) 存储效率高,同步性好。为了保证效率,Redis 将数据缓存在内存中,并周期性地把更新的数据写入磁盘或者把修改操作写入追加的记录文件中,并且在此基础上实现了主从同步。

8.4.2 Bigtable

Bigtable 是 Google 在 2004 年开始研发的一个分布式结构化数据存储系统,运用按列存储数据的方法,是一个未开源的系统。目前,已经有超过百余个项目或服务是由 Bigtable 来提供技术支持的,如 Google Analytics、Google Finance、Writely、Personalized Search 和

